

than the single descriptor approach. Subsequently, researchers began to use profiles for multiple windows. There could be two, three, four windows where the members of the family could agree on content. Sometimes, a profile was not built explicitly but rather was maintained as a collection of the instances across the known or alleged family members of the conserved region under consideration.--

Please replace the paragraph as it appears on page 13, lines 5-12 with the following rewritten paragraph:

--In step 210, the sequence threshold, K, is set. It is possible to set $K=|T|$, which is the number of sequences in the training set. In actuality, it has proven beneficial to assign a small starting value to K that is a fraction of the number of sequences in T. Experiments have shown that a starting value of $K=|T|/b$ with $b=4$ or 5 is a good choice across many data sets. Note that the smaller the value of b, the higher the redundancy of the composite descriptor will be. The selection of K also can depend on how conserved, or similar, the family members are. If the family members are well conserved, then K can be higher; if the family members are not well conserved, then K can be lower.--

Please replace the paragraph as it appears on page 26, line 23, through page 27, line 5 with the following rewritten paragraph:

--Note that the 5 hits SYV_FUGRU, GTT1_RAT, GTT1_MOUSE, SYEP_HUMAN and GTH4_MAIZE are clearly separated from the 11 top scoring sequences. They however obtained scores which were above threshold and thus are studied in more detail. In all 5 cases, one or more sizeable regions that were shared with one or more members of the PS50040 collection were discovered. The Clustal-W alignment of EF1G_XENLA and the N-terminus of SYV_FUGRU, a valyl-trna synthetase from Fugu rubripes, are shown in Table 1 below. Table 1 shows a Clustal-W alignment of EF1G_XENLA and the N-terminus of SYV_FUGRU, and this shows a strong similarity. As can be seen, the similarity among these two sequences is pretty extended and the Clustal-W score for the shown alignment equaled 462.--



Please replace Table 1 as it appears on page 28 with the following rewritten Table 1

RECEIVED
JUN 12 2003
TECH CENTER 1800/2900

--Table 1

EF1G_XENLA	(SEQ ID NO 1)	MAGGTLYTPDNWRAYKPLIAAQYSGFPIKVASSAPEFQFGVTNKTPEFLKKFPLGKVP
SYV_FUGRU_piece	(SEQ ID NO 2)	MA--TLVSP-----HLDDFRSLLALVAAEY----- ** ***, * : ..* : * **:
EF1G_XENLA	(SEQ ID NO 3)	FEGKDGFCLEFESSAIAHYVGNDELRGTTTLHQAVIQWVSFSDSHIVPPASAWVFPTLGI
SYV_FUGRU_piece	(SEQ ID NO 4)	-----C-----GNAKQ-----QSQVWQWLSFADNELTPVSCAVVFPLMGM * ** : ** ** ** ** ** ** ** ** ** **
EF1G_XENLA	(SEQ ID NO 5)	MQYNQATEQAKEGIKTVLGVLDLHQTTRTLVGERITLADITVTCSSLWLYKQVLEPSF
SYV_FUGRU_piece	(SEQ ID NO 6)	TGLDKKIQQNSRVELMRVLKVLDAQLEPRTFLVGESITLADMAVAMAVLLPFKYVLEPSD *: : : : * ***, * : ** ** ** * : : * : * **
EF1G_XENLA	(SEQ ID NO 7)	RQPFQNVTRWVFTCVNQPEFRAVLGEVKLCKMAQFQDAKKFAEMQPKKETPKKEKPAKEP
SYV_FUGRU_piece	(SEQ ID NO 8)	RNVLMNVTRWFTTCINQPEFLKVLKGISLCEKMPVTAKTSTEEAAVH-PDAAALNGPP * : * ** ** ** * : ** ** * : * : *
EF1G_XENLA	(SEQ ID NO 9)	KKEKEKKKAAPTAPAPEDDLDESEKALAAEPKSKDPYALP-KSSFIMDEFKRKYSNE
SYV_FUGRU_piece	(SEQ ID NO 10)	KTEAQLKKEAKREKLEKFOQKEMEAKKKMPVAEKKAKPEKRELGVITYDIPTPSGEK * * : ** * : : * * : * : *
EF1G_XENLA	(SEQ ID NO 11)	DTLTVALPYFW-EHFDKEGWSIWAAY-KFPEELTQAFMSCNLITGMFOR-LDKLRKTGF
SYV_FUGRU_piece	(SEQ ID NO 12)	KDVVSPLPDSYSPQYVEAAWYPWWEKQGFQKPEFGRKSIGEQNPRGIFMMCIPPPNVTGS * : * ** : : : * * : * : * : * : *
EF1G_XENLA	(SEQ ID NO 13)	ASVILFGTNNNSSISGVWV-FRQDLAFTLSED-----WQIDYESYNWRKLDGSEEC--
SYV_FUGRU_piece	(SEQ ID NO 14)	LHLGHALTNAIQDTLTRWHMRGETTLWNPGCDHAGIATQVVVEKKLMREKGTSRHDLGR * ** * : ** : : * : * : *
EF1G_XENLA	(SEQ ID NO 15)	KTLLVKEYFAWEGE-----FKNVGKPFNQG-KIFK-----
SYV_FUGRU_piece	(SEQ ID NO 16)	EKFIEEVWKKNEKGDRIYHQLKLGSSLDWDRACFTMDPKLSYAVQEAFFIRMHDEGVII : : : * : * : * : * : * : *

C4

Q7
over

Subsequently, a training set T was formed by collecting the sequences and fragments listed in the first 80 positions, arguably a very small set if one considers the diversity of the GPCR family. Essentially, slightly less than 1/10-th of the available dataset were randomly sub-selected for the purposes of building the composite descriptor. Table 5 below contains a listing of the labels of the 80 sequences in this training set. Table 5 shows the Swiss-Prot labels of the 80 sequences in the training set for the G protein-coupled receptor experiment. The labels are listed in the order they were selected and they correspond to both sequences and sequence fragments.--

Please replace the paragraph as it appears on page 43, lines 10-14 with the following rewritten paragraph:

Q8

--The three composite descriptors were used to search the collection of 19,099 ORFs that were reported for the *C. elegans* genome, by the Washington University in St. Louis, School of Medicine, Genome Sequence Center, as of June 13, 1999. In all three cases, the corresponding values of $\text{Thres}_{\text{rand}}$ that were established by searching RAND-Swiss-Prot were used.--

Please replace the paragraph as it appears on page 48, lines 4-20 with the following rewritten paragraph:

--The fragments were:

Q9

```
>Y94H6A_142.g fragment (SEQ ID NO 55)
IFDNTNDLVASLLGISSITVYRKRRKIGEE
>C16C2.1 fragment (SEQ ID NO 56)
YLSGSTRAKLAESLGLSDNQKVWFQNRRT
>F18C5.2 fragment (SEQ ID NO 57)
ISRSTAKEVATARGISEGTVSYLAMAVEK
>Y39F10A.a fragment (SEQ ID NO 58)
LSAYTISDLAKHFNVSKEILKIDIEGAEL
>Y48C3A.s fragment (SEQ ID NO 59)
NEVLNLNEVAKELNISKRRVYDVINVLEGL
```

and their respective top-scoring sequences from the training set of 70 helix-turn helix segments, blast scores, P and N values are:

#	C. elegans ORF	Top Scoring	Scor	P	N
1	Y94H6A_142.g	RPSF_BACSU	50	2.80E-06	1
2	C16C2.1	TER3_ECOLI	45	1.30E-05	1